

Active Query Selection for Crowd-Based Reinforcement Learning

Jonathan Erskine^{*†}, Taku Yamagata^{*}, Raúl Santos-Rodríguez

University of Bristol
{jonathan.erskine,taku.yamagata,enrsr}@bristol.ac.uk

Abstract

Preference-based reinforcement learning has gained prominence as a strategy for training agents in environments where the reward signal is difficult to specify or misaligned with human intent. However, its effectiveness is often limited by the high cost and low availability of reliable human input, especially in domains where expert feedback is scarce or errors are costly. To address this, we propose a novel framework that combines two complementary strategies: probabilistic crowd modelling to handle noisy, multi-annotator feedback, and active learning to prioritize feedback on the most informative agent actions. We extend the Advise algorithm to support multiple trainers, estimate their reliability online, and incorporate entropy-based query selection to guide feedback requests. We evaluate our approach in a set of environments that span both synthetic and real-world-inspired settings, including 2D games (Taxi, Pacman, Frozen Lake) and a blood glucose control task for Type 1 Diabetes using the clinically approved UVA/Padova simulator. Our preliminary results demonstrate that agents trained with feedback on uncertain trajectories exhibit faster learning in most tasks, and we outperform the baselines for the blood glucose control task.

Introduction

Preference-based reinforcement learning has become a powerful paradigm for tailoring the behaviour of autonomous agents with human trainers. By incorporating human judgments, typically in the form of pairwise comparisons indicating whether a particular action is preferred over the other, human preferences can guide exploration, improve sample efficiency, and prevent unsafe behaviours during training (Kaufmann et al. 2024). However, practical deployment remains limited by two central challenges: (1) the cost of collecting high-quality feedback from crowds of humans with mixed skills and motives and (2) the limited availability of reliable human trainers, particularly in domains where expertise is scarce or errors carry high risk (e.g., healthcare,

robotics).

One approach to mitigate the scarcity of expert feedback is to aggregate input from multiple non-expert trainers, a technique broadly studied in the learning from crowds literature. These methods aim to infer both the true label and trainer reliability by modelling the noise inherent in human inputs. However, while crowd-based strategies can scale feedback collection (Buecheler et al. 2010), not all trainers are equally reliable.

A second, complementary strategy to reduce annotation burden is active learning. In supervised settings, active learning algorithms prioritize annotation of the most uncertain or informative data points to accelerate training with fewer labels. In RLHF, however, relatively little work has been done to adapt such ideas to interactive feedback-driven learning. Yet, not all agent behaviours are equally valuable to label; querying human feedback only on high-uncertainty actions could yield greater improvements in policy quality for the same feedback budget.

Many recent works collect feedback as pairwise comparisons of agent actions or trajectories (Abdelkareem, Shehata, and Karray 2022) and tend to convert these comparisons into a reward signal (Chakraborty et al. 2024) and/or generate models which can learn to predict these preferences (Frick et al. 2024). In our experiments we instead prescribe to the practice of collecting feedback directly on the policy as an additional signal to environmental rewards, as demonstrated by Griffith et al. (2013).

In this work, we propose a novel feedback aggregation framework that combines active learning, learning from crowds, and policy feedback: probabilistic crowd modelling to handle noisy, multi-trainer feedback, and entropy-based active selection of state-action pairs to maximize the utility of human input. Building on the framework provided in (Griffith et al. 2013), our method aims to:

1. Handle feedback from multiple trainers of varying reliability.
2. Infer trainer consistency levels in an online setting.
3. Actively request feedback only for agent behaviours deemed highly uncertain.

Figure 1 illustrates our proposed method. To evaluate our approach, we simulate human feedback using a pre-trained

^{*}These authors contributed equally.

[†]Jonathan Erskine’s PhD research is jointly funded by UK Research and Innovation (UKRI) and Thales Training & Simulation Ltd. through the UKRI Centre for Doctoral Training in Interactive Artificial Intelligence under grant EP/S022937/1. This work was partially funded by the UKRI Turing AI Fellowship EP/V024817/1.

oracle and assess performance across a set of diverse domains, including classic 2D environments (Taxi, PACMAN, Frozen Lake) and a blood-glucose monitoring and control task for Type 1 Diabetes using the clinically approved UVA/-Padova simulator (Man et al. 2014). Experimental results show that our method significantly improves learning speed and feedback efficiency, particularly in deterministic settings where uncertainty is easier to quantify and resolve.

Related Work

Learning From Crowds

Learning from crowds is a well-established research area that addresses the challenge of aggregating noisy, unreliable, or biased labels provided by multiple trainers, often with varying expertise. A naive approach treats all trainers as equally reliable. However, performance can be improved by adjusting for trainer reliability (Snow et al. 2008).

One of the earliest attempts to compute trainer reliability explicitly was a probabilistic model that jointly estimates the true labels and the error rates of each trainer via Expectation-Maximization (EM) (Dawid and Skene 1979). This foundational model has since inspired a range of probabilistic and Bayesian methods (Whitehill et al. 2009; Raykar et al. 2010), which learn trainer-specific confusion matrices and incorporate priors over worker accuracy.

Beyond probabilistic modelling, recent work leverages deep learning to model crowd annotations at scale. Rodrigues and Pereira (2018) proposed learning with crowds using Gaussian processes and variational inference, while more recent neural architectures integrate crowd modelling directly into end-to-end training (Guan, Houlsby, and Cesa-Bianchi 2018; Chu, Ma, and Wang 2020), enabling better generalization in image classification and NLP tasks. In the reinforcement learning domain, Siththaranjan, Laidlaw, and Hadfield-Menell (2024) address the challenge of collecting preferences from humans with diverse preferences by modelling preference distributions and identifying hidden contexts and improving adversarial robustness.

Active Learning

Active learning (AL) is a supervised learning framework that aims to reduce annotation cost by selecting only the most informative data points for labelling. This selection is typically guided by an acquisition function that quantifies a model’s uncertainty or expected gain from receiving a label. The foundational idea is that the learner queries those instances for which its current predictions are least confident. This approach, known as *uncertainty sampling*, was introduced by Lewis and Gale (1994). Subsequent strategies have proposed entropy-based selection (using Shannon information), margin sampling (difference between top prediction confidences), and expected model change (querying points likely to cause the greatest update to the model) (Settles 2009). These techniques have proven effective across a wide range of domains, including natural language processing, computer vision, and structured prediction.

While traditionally applied to classification tasks, active learning has recently seen broader use in more complex

learning settings. Notably, several recent works have formulated active learning itself as a sequential decision-making process using reinforcement learning to learn acquisition policies (Fang, Li, and Cohn 2017; Ebert, Fritz, and Schiele 2012). Other works in the RLHF domain adopt reward model training as the target task, where feedback is derived from human preference annotations and refined through active sampling (Ouyang et al. 2022).

In our approach, we use active learning only at the feedback selection stage: after each episode, the agent identifies the most uncertain state-action pairs and queries for feedback on only those. This selective feedback helps improve sample efficiency and avoids unnecessary burden on human trainers.

Reinforcement Learning with Human Feedback

Rather than relying solely on environment-provided rewards, Reinforcement learning with human feedback (RLHF) agents are guided by human evaluators through signals such as scalar feedback, preferences, or demonstrations. Early work in this area includes TAMER (Knox and Stone 2012), which proposed learning from real-time human-provided scalar feedback. Deep TAMER (Warnell et al. 2018) extended this framework to high-dimensional continuous domains by using deep networks to model human responses. DQN-TAMER (Arakawa et al. 2018) further integrated TAMER-style feedback with standard deep RL objectives. Other methods such as Advise (Griffith et al. 2013) introduced the idea of policy shaping, directly modifying the agent’s action probabilities based on binary human input. An influential shift came with preference-based methods (Christiano et al. 2017), where agents learn a reward model from human comparisons between trajectory segments, rather than receiving direct rewards or feedback. This line of work is exemplified by PEBBLE (Lee, Smith, and Abbeel 2021), which first trains agents to explore a wide range of states and then elicits human preferences to learn a reward function. Recent work such as DUO (Feng et al. 2025) advances this idea by incorporating active query selection into RLHF. DUO selects trajectory pairs based on informativeness, diversity, and proximity to the agent’s current policy (on-policy sampling), thereby improving the sample efficiency of reward learning. Similarly Chhan, Novoseller, and Lawhern (2025) use methods derived from unsupervised ensemble learning to effectively aggregate user preference feedback across diverse crowds to learn RL policies. However, these approaches typically assume that no environment reward is available and thus focus entirely on constructing an accurate reward model from human input.

By contrast, our approach complements the environment reward with targeted human feedback. We assume the agent already receives a reward signal from the environment, and our goal is to augment this with feedback on only the most uncertain state-action pairs. Unlike most RLHF approaches that model uncertainty at the level of the reward function, our method uses entropy derived from the value function to guide query selection.

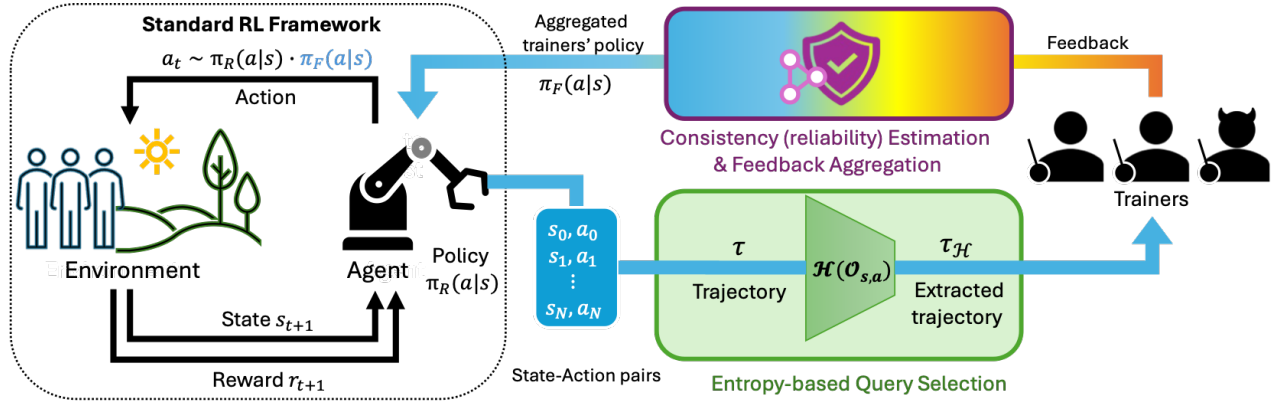


Figure 1: Standard reinforcement learning training loop and our extended method. We collect a trajectory τ of state–action pairs and extract a subset $\tau_{\mathcal{H}} \subseteq \tau$ consisting of the top- n pairs ranked by entropy $\mathcal{H}(\mathcal{O}_{s,a})$. This subset is labelled by a crowd of trainers with diverse skill levels (some might be adversarial). We estimate each trainer’s reliability based on their feedback and weight their input accordingly when aggregating their feedback. Finally, the aggregated results are merged with the underlying RL policy $\pi_R(a | s)$, allowing us to sample an action to execute in the environment. This approach improves the efficiency of human feedback while accounting for varying trainer reliability.

Preliminary

Reinforcement Learning

Reinforcement learning (RL) is a framework for any learning process that involves sequentially interacting with an environment to achieve a certain objective (Sutton and Barto 1998). The learner is called the *agent*. It performs an action a_t in the *environment*, observes its consequences s_{t+1} , and receives a reward r_{t+1} (or a cost) signal – a numerical assessment of the current situation as shown in the standard RL framework in Fig. 1.

The agent basically learns a mapping from the state to the action that maximises the total amount of reward it receives over the long run. The mapping is called *policy* denoted as $\pi_t(a | s)$ that indicates the probability of $a_t = a$ when the state is $s_t = s$.

Advise Algorithm

In order to incorporate feedback into the RL algorithm, we build upon the *Advise* algorithm (Griffith et al. 2013) and thus provide a brief description of the approach. *Advise* assumes binary feedback from a trainer that returns either ‘right’ or ‘wrong’ for a particular agent’s choice of action. The feedback is accumulated for each state-action pair separately, and it is used to derive a trainer’s policy, denoted by $\pi_F(a|s)$, which is then used to modify the agent’s policy. Additionally, C is *consistency level*, defined as the probability that the trainer gives the right (consistent) feedback, and assuming a binomial distribution, the authors propose the following trainer policy.

$$\pi_F(a | s) \propto C^{\Delta(s,a)}(1 - C)^{-\Delta(s,a)} \quad (1)$$

where $\Delta(s, a)$ is the difference between the number of positive and negative feedback from the trainer.

The policy of the trainer is combined with $\pi_R(s, a)$ (policy from the underlying RL algorithm) by multiplying them

together so that the final policy becomes as

$$\pi(a | s) \propto \pi_F(a | s) \cdot \pi_R(a | s). \quad (2)$$

This formulation enables direct shaping of the agent’s behaviour through human feedback, rather than modifying the reward function. Such an approach is often considered more effective, as it allows the feedback to provide immediate influence on the policy itself.

Method

While *Advise* provides a mechanism for incorporating human feedback, it exhibits three key limitations. First, it requires prior knowledge of the trainer’s consistency level (C), which may be unknown or difficult to estimate in many practical applications. Second, the algorithm assumes a single trainer, limiting its scalability in settings involving multiple human supervisors or crowd-sourced feedback. Third, *Advise* typically requires a large amount of feedback before producing noticeable improvements. While this is not necessarily a flaw of the algorithm itself, it presents a significant barrier in real-world scenarios where human feedback is costly and limited.

In this section, we propose a method to address all of the limitations above. Our approach supports feedback from multiple trainers and estimates each trainer’s reliability (i.e., consistency level) in an online manner as feedback is collected (feedback from crowds). Additionally, the algorithm actively queries for feedback in situations where it is expected to be most informative (entropy-based active feedback). By estimating trainer reliability on the fly, the system leverages high-quality feedback while mitigating the impact of unreliable or adversarial trainers, making it robust to feedback of varying quality. Furthermore, by strategically selecting when to request feedback, the algorithm enhances the efficiency of human supervision, reducing the overall burden on the trainers.

Feedback from Crowds

Learning from Crowds is a well-established approach in supervised learning where labels are collected from a group of non-experts, reducing the burden of acquiring expert annotations. Typically, a probabilistic model iteratively estimates both the true labels and the reliability of each trainer to handle low-quality (noisy) annotators. We extend this framework to the reinforcement learning (RL) setting with introducing a variational inference (VI) approach. This method iteratively estimates the distributions over both the optimal actions for given states and the reliability of individual trainers. The use of VI provides two key benefits: it allows the integration of prior knowledge, thereby improving the sample efficiency of the inference process, and it models uncertainty through parameter distributions, enhancing robustness in scenarios with sparse feedback.

Our VI-based approach infers the posterior distributions of trainer reliability (consistency level, C_l) and the optimality of state-action pairs ($\mathcal{O}_{s,a}$) by maximizing the likelihood of human feedback observations. These observations are represented by $h_{l,s,a}^+$ and $h_{l,s,a}^-$, denoting the counts of positive and negative feedback, respectively, given by trainer l for state s and action a . The remainder of this section provides a high-level overview of our method. We refer the reader to Appendix A for additional details.

First, we define the posteriors using two variational factors $q(\mathbf{O}_s)$ and $q(C_l)$, as follows:

$$P(\mathbf{O}, \mathbf{C} | \mathbf{h}^+, \mathbf{h}^-) = \left(\prod_s q(\mathbf{O}_s) \right) \left(\prod_l q(C_l) \right), \quad (3)$$

where $\mathbf{O}, \mathbf{C}, \mathbf{h}^+$ and \mathbf{h}^- denote collections of $\mathcal{O}_{s,a}, C_l, h_{l,s,a}^+$ and $h_{l,s,a}^-$, respectively for all trainers, states and actions. \mathbf{O}_s represents the set of $\mathcal{O}_{s,a}$ values for all actions in a given state s .

The VI approach updates the estimates of $q(\mathbf{O}_s)$ and $q(C_l)$ iteratively by maximising the evidence lower bound (ELBO) until convergence, after which we employ the *posterior sampling* (Strens 2000; Thompson 1933) by treating the estimated posterior $q(\mathbf{O}_s)$ as the policy:

$$\pi(a|s) = q(\mathbf{O}_s = \mathbf{e}_a), \quad (4)$$

where \mathbf{e}_a is a one-hot encoding vector. We take the underlying RL policy as the prior of $q(\mathbf{O}_s)$. Consequently, it is naturally incorporated into the the posterior and the final policy is based on both interactions with environment and human feedback.

Entropy-based Active Feedback

Our method builds on an entropy-based measure of the posterior distribution over optimality variables for state-action pairs. Specifically, we compute the entropy of these variables along the experienced trajectory and prioritise feedback requests for the action pairs showing high entropy.

The posterior of the optimality variable $\mathcal{O}_{s,a}$ is conditioned on two types of observations: human feedback and interactions with the environment. Formally, the posterior is denoted as $P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-, \tau)$, where the \mathbf{h}^+ and \mathbf{h}^-

represent positive and negative feedback, respectively and τ is a set of trajectories generated from interacting with the environment. Each trajectory τ consists of a sequence of state s_t , action a_t and reward r_t at each time step t e.g. $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \dots\}$.

Assuming conditional independence between human feedback $\mathbf{h}^+, \mathbf{h}^-$ and the trajectories τ given the optimality variable, the posterior can be factorised as follows (a detailed derivation is provided in Appendix B):

$$P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-, \tau) \propto P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-) P(\mathcal{O}_{s,a} | \tau) / P(\mathcal{O}_{s,a}). \quad (5)$$

Here $P(\mathcal{O}_{s,a})$ denotes the prior over the optimality. We assume a uniform prior, such that $P(\mathcal{O}_{s,a} = 1) = \frac{1}{N_a}$, where N_a is a number of actions in a state. The two posterior components – one conditioned on human feedback and the other on interactions with the environment – can be computed independently. These components are then combined multiplicatively, adjusted by the inverse of the prior, to yield the joint posterior. This structure facilitates scalable inference and enables integration of heterogeneous information sources.

In the following sections, we describe the procedures used to derive each of these posterior terms.

The posterior given the human feedback $P(\mathcal{O}_{s,a} | \mathbf{H})$
We build upon the results of *Advise* (Griffith et al. 2013), using its formulation to derive the posterior given human feedback, as expressed in the following equation:

$$P(\mathcal{O}_{s,a} | \mathbf{H}) \propto C^{\Delta(s,a)} (1 - C)^{-\Delta(s,j)}, \quad (6)$$

where $\Delta_{s,a}$ represents the difference between the number of positive and negative feedback on state s and action a .

The posterior given the trajectories $P(\mathcal{O}_{s,a} | \tau)$
We consider a baseline RL algorithm that estimates the state-action value function $Q(s, a)$ from the trajectories. In this context, an action is considered optimal if it yields the highest value amongst all available actions in a given state. Therefore, we define the posterior distribution of the optimality variable as:

$$P(\mathcal{O}_{s,a} | \tau) = P(Q(s, a) > Q(s, a') \text{ for all } a' \neq a). \quad (7)$$

To make this computation tractable, we assume conditional independence of $Q(s, a')$ for all $a' \neq a$ given $Q(s, a)$, and factorise each inequality. While this assumption is strong, it is reasonable if the actions are sufficiently distinct from other. This factorisation simplifies the posterior computation, as it allows us to express it as a product of independent terms.

$$\begin{aligned} P(\mathcal{O}_{s,a} | \tau) &= \int P(Q(s, a) = X) \cdot \\ &\quad P(X > Q(s, a') \text{ for all } a' \neq a | X) dX \\ &= \mathbb{E}_{X \sim P(Q(s,a))} \left[\prod_{a' \neq a} P(X > Q(s, a')) \right] \\ &= \mathbb{E}_{X \sim P(Q(s,a))} \left[\prod_{a' \neq a} F_{s,a'}(X) \right]. \end{aligned} \quad (8)$$

Here, $F_{s,a'}(X)$ denotes the cumulative distribution function (CDF) of $P(Q(s, a'))$.

Next, we assume that $P(Q(s, a))$ follows Gaussian distribution with the mean $\hat{Q}(s, a)$ and the standard deviation $\sigma_{base}/\sqrt{N_{s,a}}$, where $\hat{Q}(s, a)$ is the estimated value function, σ_{base} is a baseline standard deviation (hyper-parameter) and $N_{s,a}$ is a number of visits to the state-action pair.

We compute the expectation in Eq. 8 using Monte Carlo estimation. Specifically, we apply the inverse transform method with stratified sampling. First, we construct a sequence of M evenly spaced probability values:

$$p_i = \frac{1}{2M} + \frac{i}{M} \text{ for } i = 0, 1, 2, \dots, M-1. \quad (9)$$

Then, we obtain the samples of $Q(s, a)$ by applying the inverse CDF:

$$x_i = F_{s,a}^{-1}(p_i) \text{ for } i = 0, 1, 2, \dots, M-1. \quad (10)$$

Finally, we approximate the posterior using:

$$P(\mathcal{O}_{s,a}|\tau) \simeq \frac{1}{M} \sum_{i=0}^{M-1} \prod_{a' \neq a} F_{s,a'}(x_i). \quad (11)$$

Since we assume a Gaussian distribution for $Q(s, a)$, both $F_{s,a}^{-1}(p_i)$ and $F_{s,a'}(x_i)$ can be computed analytically. In addition, the Monte Carlo estimation process lends itself well to parallel (vectorised) computation, which improves efficiency and reduces runtime.

One-vs-All Entropy We now have the posterior distribution $P(\mathcal{O}_{s,a}|\mathbf{h}^+, \mathbf{h}^-, \tau)$ for each state-action pair, and we are ready to compute its entropy. However, the conventional Shannon entropy is not well-suited in this context. The variable $\mathcal{O}_{s,a}$ is binary where $\mathcal{O}_{s,a} = 1$ indicates the action a is optimal at state s , and $\mathcal{O}_{s,a} = 0$ implies that the other action(s) is optimal. Consequently, we require an entropy measure \mathcal{H} that satisfies the following criteria (Requirement IV is desirable but not essential):

- I. $\mathcal{H}(\mathcal{O}_{s,a}) \rightarrow 0$ when $P(\mathcal{O}_{s,a} = 1) \rightarrow 0$.
- II. $\mathcal{H}(\mathcal{O}_{s,a}) \rightarrow 0$ when $P(\mathcal{O}_{s,a} = 1) \rightarrow 1$.
- III. $\mathcal{H}(\mathcal{O}_{s,a})$ reaches its max. when $P(\mathcal{O}_{s,a} = 1) = 1/N_a$.
- IV. $\mathcal{H}(\mathcal{O}_{s,a})$ reduces to the Shannon entropy when $N_a = 2$.

To meet these requirements, we propose the One-vs-All (OvA) entropy. The OvA entropy is defined for a binary variable ($\mathcal{O}_{s,a}$) with Bernoulli distribution (parameter p) and it is computed in two steps.

1. Re-normalise Bernoulli parameter p using: $p' = \frac{p}{p+(1-p)/(N_a-1)}$.
2. Compute Shannon entropy for Bernoulli distribution with the re-normalised parameter p' .

The resulting OvA entropy \mathcal{H}_{OvA} is given by:

$$\mathcal{H}_{OvA}(\mathcal{O}_{s,a}) = -p' \log_2(p') - (1-p') \log_2(1-p'), \quad (12)$$

where $p' = \frac{p}{p+(1-p)/(N_a-1)}$ and $p = P(\mathcal{O}_{s,a} = 1|\mathbf{h}^+, \mathbf{h}^-, \tau)$. The re-normalisation scales

down the non-optimal probability $1-p$ by the number of other actions N_a-1 . As the result, the OvA entropy \mathcal{H}_{OvA} gives the maximum value when $p = 1/N_a$. This form of entropy is related in nature to the family of asymmetric entropies (Guermaz, Chaabane, and Hammami 2018; Santos-Rodríguez et al. 2009).

Summary The entropy-based active feedback method proceeds as follows: first, the agent collects trajectories (state-action pairs) through interaction with the environment. After completing an episode (or after a fixed number of episodes), the method computes the posteriors of the optimality variable for all state-action pairs encountered in the trajectories, using Eq. 5, Eq. 6, and Eq. 11. Next, the OvA entropy is calculated for each state-action pair. Finally, the method identifies the state-action pairs with the top highest OvA entropies and requests human feedback for these pairs.

This method strategically prioritises querying human feedback for state-action pairs exhibiting high uncertainty with respect to their optimality. As a result, the system maximises the information gain from human feedback by targeting ambiguous or uncertain decisions. Moreover, by limiting the scope to state-action pairs encountered within actual agent trajectories, the approach concentrates learning on decision points that are relevant and likely to recur during future interactions with the environment.

Experiments

We conduct experiments across a diverse set of environments that operate on a gridworld domain; these include our own PACMAN-inspired gridworld as well as Taxi and FrozenLake from the Farama Gymnasium suite (Towers et al. 2024). We then test our method on a real-world task; blood glucose control for Type 1 Diabetes using the clinically validated UVA/Padova simulator. Across all environments, we compare our entropy-based active querying and crowd feedback aggregation approach against (1) the baseline approach (no feedback) and (2) random active learning (feedback randomly sampled from trajectories). We train an oracle policy using standard reinforcement learning on the target environment as a proxy for an expert human. We modify $C_t = [0.9, 0.8, 0.6, 0.3]$ across 4 trainers to simulate human trainers of varying quality.

Gridworld Domains

We first implemented a simplified abstraction of the classic PACMAN game, using a small (5x5) grid and a reduced state space consisting of two pellets and a single ghost. This controlled setup allows us to test agent learning in environments that require sequential planning, without the full complexity of the original game.

We validated our **consistency level estimation** method in this PACMAN environment. After confirming the effectiveness, we expanded our evaluation to additional gridworld domains with varying dynamics and challenges:

- **Taxi environment:** complementary evaluation featuring random starting positions and destinations and a larger map. Default parameters from the gymnasium library were used without modification

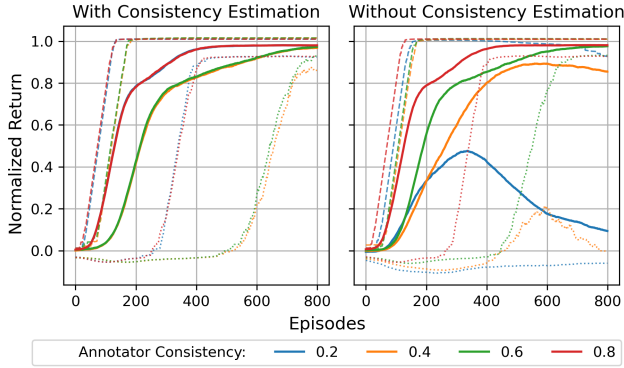


Figure 2: Evaluation results with and without the consistency level estimation. The agent without the estimation assumes the consistency level is 0.8 fixed, while varying the true value over 0.8, 0.6, 0.4 and 0.2. Dashed and dotted lines correspond to the 95% upper and 5% lower confidence bounds, respectively.

- **Customised Frozen Lake environments:** designed to investigate the relationship between task difficulty and performance of our active learning method.

For each environment, we run 50 trials of 1000 episodes per agent and take the average return to account for training noise. Our learning rate is set to $\alpha = 0.05$ and discount factor $\gamma = 0.9$. Detailed environment descriptions are provided in the technical appendix. Our code and experiments will be made publicly available in the camera-ready version.

Results We validated our **consistency level estimation** method in the PACMAN environment by comparing agents with and without this feature without using our active feedback. In the non-estimating condition, we fix the consistency level at 0.8, while varying the true value across 0.8, 0.6, 0.4, and 0.2. Figure 2 shows that agents without estimation perform poorly when the assumed and true consistency differ. Agents with the ability to estimate consistency perform well across all values. The best overall performance occurs when the true level is 0.8 (which is matched with the assumed level) and estimation is disabled.

These results demonstrate that consistency level estimation enhances robustness when the agent receives feedback from trainers with unknown reliability. However, if the trainer’s consistency level is known, directly using the known value yields better performance than estimation.

We now extend our evaluation to our full suite of environments incorporating both consistency level estimation and active feedback mechanisms. Figure 3 shows learning curves for PACMAN and Taxi. In PACMAN, our entropy-based active learning method (AL-Entropy) outperforms both the baseline and random sampling (AL-Random). In contrast, results on Taxi are less conclusive. We hypothesise that this discrepancy is due to the structural differences between the environments: Taxi presents a relatively unconstrained space with many viable paths to success, while PACMAN imposes stricter conditions for winning; particularly due to board sec-

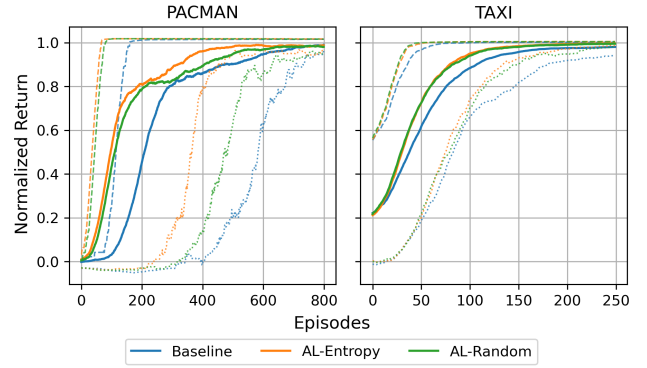


Figure 3: Evaluation results for PACMAN and Taxi grid-world environments. Agents are trained with standard Q-Learning (Baseline) and active learning with random sampling (AL-Random) and entropy-based sampling (AL-Entropy). Dashed and dotted lines correspond to the 95% upper and 5% lower confidence bounds, respectively.

Environment	Baseline	AL-Random	AL-Entropy
PACMAN	546.91	632.02	660.77
Taxi	201.33	209.19	209.75
Frozen Lake (0)	0.00	32.88	34.76
Frozen Lake (1)	0.24	38.30	38.59
Frozen Lake (2)	-0.01	31.38	35.14
Frozen Lake (3)	0.00	19.12	26.40

Table 1: Area Under the Curve (AUC) comparison across gridworld environments and methods. Higher values indicate better performance.

tions where the goal (a pellet) is difficult to reach without navigating around a ghost.

To explore this hypothesis, we test four variants of the Frozen Lake environment: one with a randomly generated map offering multiple solution paths, and three increasingly constrained maps where successful navigation depends on passing through one or more mandatory “gates”.

Figure 4 illustrates the learning curves for entropy-based sampling vs. random sampling in these environments. Both methods perform similarly in the open, unconstrained setting. There is no discernable improvement from random sampling in the easy environment (1). However, as the environment becomes more constrained, entropy-based sampling increasingly outperforms random sampling. Table 1 shows the difference in area under the curve (AUC) for each of our experiments. A larger AUC implies faster learning.

Type 1 Diabetes BGL Control

We evaluate our approach in a clinically relevant healthcare application: the control of blood glucose level (BGL) in individuals with type 1 diabetes. The objective of this task is to maintain BGL within a healthy range by administering insulin at the appropriate dosage and timing. The agent mon-

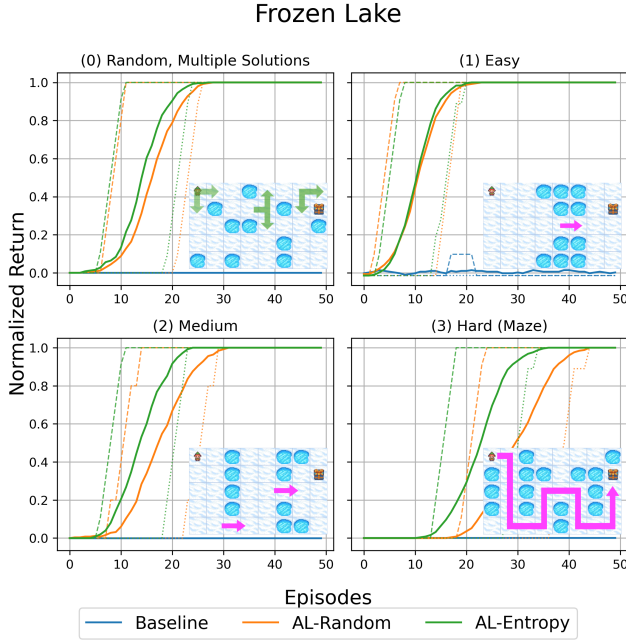


Figure 4: Evaluation results for the Frozen Lake grid-world environment. Agents are trained with standard Q-Learning (Baseline) and active learning with random sampling (AL-Random) and entropy-based sampling (AL-Entropy). Dashed and dotted lines correspond to the 95% upper and 5% lower confidence bounds, respectively.

itors BGL through continuous glucose monitor (CGM) and receives information about carbohydrate intake (i.e., meals). Based on these observations, it determines both the timing and amount of insulin to be delivered.

For this evaluation, we employ the UVA/Padova type 1 diabetes simulator (Man et al. 2014), the first computational model approved by the U.S. Food and Drug Administration (FDA) as a substitute for preclinical trials of certain insulin treatments. We use an open-source implementation of this simulator (Xie. 2018), which includes physiological profiles of 30 distinct virtual individuals. The simulator models the physiological response of BGL to insulin administration, allowing us to systematically assess the agent’s performance in a realistic and high-fidelity simulation environment.

Results We evaluate three agent models for this task: standard Q-learning (Baseline), Q-learning with feedback with uniformly random sampling (AL-Random), and Q-learning with entropy-based active feedback (AL-Entropy). We test these models on three virtual individuals provided by the UVA/Padova simulator. Each model is trained on 2,000 episodes with five random seeds, and performance is measured by the percentage of time that BGL remain within the target range of 70–180 mg/dL during the training period.

The results, presented in Table 2, show that the entropy-based active feedback model consistently outperforms the other baselines across subjects, confirming the benefits of the approach on a challenging safety-critical task.

Individual	Baseline	AL-Random	AL-Entropy
child003	36.7 \pm 1.4	49.5 \pm 1.6	54.4 \pm 1.1
adolescent003	45.7 \pm 3.6	69.6 \pm 1.9	75.6 \pm 0.9
adult003	60.6 \pm 4.3	80.4 \pm 3.7	88.6 \pm 1.7

Table 2: % of time in the target BGL range (70-180 mg/dL) across entire training episodes (2000) with the standard deviation. The best results amongst the three models are highlighted in bold face.

Discussion

Our gridworld experiments enabled us to test environments with different dynamics. We propose that in environments with many good solutions, the uncertainty in any given state might be spread across multiple acceptable actions; the entropy for each action might be relatively high, but not because one action is definitively better than another. In such a scenario, randomly sampling might be just as effective, or even more efficient than calculating entropy. Entropy-based active learning surpasses random sampling when the problem becomes more constrained; our superior performance on the blood glucose control task aligns with our conclusion that active learning is well-suited to these highly-constrained tasks.

Limitations The additional steps in our pipeline come at a cost; whether it is the computational effort required to process and downsample trajectories, or the time required by a (human or synthetic) teacher to annotate these trajectories, we must consider the trade-off vs. the cost of additional training time. Additionally, we assume that we have readily available expert and non-expert teachers.

Real humans could potentially lead to injecting biases into our models; any real-world application would require testing against adversarial teachers. There is also a known trade-off between the sample efficiency of active learning and the robustness of the resulting policy. While our method improves learning speed, we have yet to establish a methodology for evaluating robustness in this context.

Finally, while we have validated our method on a real-world control task, its implementation relies on tabular Q-learning; a simple and limited learning algorithm that does not scale well to more complex or high-dimensional environments.

Future Work To assess the broader applicability of our approach, future work will explore its integration with function approximation methods such as Deep Q-Networks (DQNs).

Beyond scalability, agent stability warrants further investigation alongside consideration of recent methods to enhance robustness (Yamagata and Santos-Rodriguez 2024).

Finally, we aim to implement human data collection to better understand the relationship between annotation effort and performance improvement.

References

- Abdelkareem, Y.; Shehata, S.; and Karray, F. 2022. Advances in Preference-based Reinforcement Learning: A Review. In *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2527–2532. IEEE.
- Arakawa, R.; Kobayashi, S.; Unno, Y.; Tsuboi, Y.; and Maeda, S.-i. 2018. Dqn-tamer: Human-in-the-loop reinforcement learning with intractable feedback. *arXiv preprint arXiv:1810.11748*.
- Buecheler, T.; Sieg, J.; Füchslin, R.; and Pfeifer, R. 2010. Crowdsourcing, Open Innovation and Collective Intelligence in the Scientific Method: A Research Agenda and Operational Framework. *Artificial Life XII – Twelfth International Conference on the Synthesis and Simulation of Living Systems, Odense, Denmark*, 679–686.
- Chakraborty, S.; Qiu, J.; Yuan, H.; Koppel, A.; Manocha, D.; Huang, F.; Bedi, A.; and Wang, M. 2024. MaxMin-RLHF: Alignment with Diverse Human Preferences. In *Forty-first International Conference on Machine Learning*.
- Chhan, D.; Novoseller, E.; and Lawhern, V. J. 2025. Crowd-PrefRL: Preference-Based Reward Learning from Crowds. *arXiv:2401.10941*.
- Christiano, P. F.; Leike, J.; Brown, T.; Martic, M.; Legg, S.; and Amodei, D. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Chu, Z.; Ma, J.; and Wang, H. 2020. Learning from Crowds by Modeling Common Confusions. *CoRR*, abs/2012.13052.
- Dawid, A. P.; and Skene, A. M. 1979. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 28(1): 20–28.
- Ebert, S.; Fritz, M.; and Schiele, B. 2012. RALF: A reinforced active learning formulation for object class recognition. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 3626–3633.
- Fang, M.; Li, Y.; and Cohn, T. 2017. Learning how to Active Learn: A Deep Reinforcement Learning Approach. *arXiv:1708.02383*.
- Feng, X.; Jiang, Z.; Kaufmann, T.; Xu, P.; Hüllermeier, E.; Weng, P.; and Zhu, Y. 2025. DUO: Diverse, Uncertain, On-Policy Query Generation and Selection for Reinforcement Learning from Human Feedback. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 16604–16612.
- Frick, E.; Li, T.; Chen, C.; Chiang, W.-L.; Angelopoulos, A. N.; Jiao, J.; Zhu, B.; Gonzalez, J. E.; and Stolica, I. 2024. How to Evaluate Reward Models for RLHF. *arXiv:2410.14872*.
- Goyal, A.; Islam, R.; Strouse, D.; Ahmed, Z.; Botvinick, M.; Larochelle, H.; Bengio, Y.; and Levine, S. 2023. InfoBot: Transfer and Exploration via the Information Bottleneck. *arXiv:1901.10902*.
- Griffith, S.; Subramanian, K.; Scholz, J.; Isbell, C. L.; and Thomaz, A. L. 2013. Policy Shaping: Integrating Human Feedback with Reinforcement Learning. In Burges, C.; Bottou, L.; Welling, M.; Ghahramani, Z.; and Weinberger, K., eds., *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
- Guan, M.; Hounsby, N.; and Cesa-Bianchi, N. 2018. Who said what: Modeling individual annotator reliability and annotator bias. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Guerhazi, R.; Chaabane, I.; and Hammami, M. 2018. AE-CID: Asymmetric entropy for classifying imbalanced data. *Information Sciences*, 467: 373–397.
- Kaufmann, T.; Weng, P.; Bengs, V.; and Hüllermeier, E. 2024. A Survey of Reinforcement Learning from Human Feedback. *arXiv:2312.14925*.
- Knox, W. B.; and Stone, P. 2012. Reinforcement learning from simultaneous human and MDP reward. 475–482. *International Foundation for Autonomous Agents and Multiagent Systems*. ISBN 0981738117.
- Kumar, A.; Hong, J.; Singh, A.; and Levine, S. 2022. When Should We Prefer Offline Reinforcement Learning Over Behavioral Cloning? *arXiv:2204.05618*.
- Lee, K.; Smith, L.; and Abbeel, P. 2021. Pebble: Feedback-efficient interactive reinforcement learning via relabeling experience and unsupervised pre-training. *arXiv preprint arXiv:2106.05091*.
- Lewis, D. D.; and Gale, W. A. 1994. A Sequential Algorithm for Training Text Classifiers. *arXiv:cmp-lg/9407020*.
- Liu, H.; Zhuge, M.; Li, B.; Wang, Y.; Faccio, F.; Ghanem, B.; and Schmidhuber, J. 2023. Learning to Identify Critical States for Reinforcement Learning from Videos. *arXiv:2308.07795*.
- Man, C. D.; Micheletto, F.; Lv, D.; Breton, M.; Kovatchev, B.; and Cobelli, C. 2014. The UVA/PADOVA type 1 diabetes simulator: New features. *Journal of Diabetes Science and Technology*, 8: 26–34.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C. L.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; Schulman, J.; Hilton, J.; Kelton, F.; Miller, L.; Simens, M.; Askell, A.; Welinder, P.; Christiano, P.; Leike, J.; and Lowe, R. 2022. Training language models to follow instructions with human feedback. *arXiv:2203.02155*.
- Raykar, V. C.; Yu, S.; Zhao, L. H.; Valadez, G. H.; Florin, C.; Bogoni, L.; and Moy, L. 2010. Learning from crowds. *Journal of Machine Learning Research*, 11: 1297–1322.
- Rodrigues, F. C.; and Pereira, F. 2018. Deep learning from crowds. In *AAAI Conference on Artificial Intelligence*, volume 32.
- Santos-Rodríguez, R.; Guerrero-Curieses, A.; Alaiz-Rodríguez, R.; and Cid-Sueiro, J. 2009. Cost-sensitive learning based on Bregman divergences. *Machine Learning*, 76(2-3): 271–285.
- Settles, B. 2009. Active Learning Literature Survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison.

- Siththaranjan, A.; Laidlaw, C.; and Hadfield-Menell, D. 2024. Distributional Preference Learning: Understanding and Accounting for Hidden Context in RLHF. arXiv:2312.08358.
- Snow, R.; O’Connor, B.; Jurafsky, D.; and Ng, A. Y. 2008. Cheap and fast—but is it good? Evaluating non-expert annotations for natural language tasks. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, 254–263. Association for Computational Linguistics.
- Strens, M. 2000. A Bayesian framework for reinforcement learning. volume 2000, 943–950. Posterior Sampling RL.
- Sutton, R. S.; and Barto, A. G. 1998. *Reinforcement Learning*. Cambridge, MA: The MIT Press.
- Thompson, W. R. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4): 285–294.
- Towers, M.; Kwiatkowski, A.; Terry, J.; Balis, J. U.; Cola, G. D.; Deleu, T.; Goulão, M.; Kallinteris, A.; Krimmel, M.; KG, A.; Perez-Vicente, R.; Pierré, A.; Schulhoff, S.; Tai, J. J.; Tan, H.; and Younis, O. G. 2024. Gymnasium: A Standard Interface for Reinforcement Learning Environments. arXiv:2407.17032.
- Warnell, G.; Waytowich, N.; Lawhern, V.; and Stone, P. 2018. Deep tamer: Interactive agent shaping in high-dimensional state spaces. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Whitehill, J.; Ruvolo, P.; Wu, T.; Bergsma, J.; and Movellan, J. R. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *Advances in Neural Information Processing Systems*, volume 22.
- Xie., J. 2018. Simglucose v0.2.1 <https://github.com/jxx123/simglucose>.
- Yamagata, T.; and Santos-Rodriguez, R. 2024. Safe and Robust Reinforcement Learning: Principles and Practice. arXiv:2403.18539.

Appendix A: Variational Inference Approach

We introduce our VI algorithm for the consistency level estimation. The VI algorithm estimates the distribution of the consistency level. Hence, it captures the uncertainty of the estimation, and we expect it to be more robust against potential estimation errors, especially when we do not receive enough feedback from the trainer. Also, VI could incorporate prior information, and it potentially dramatically improves the estimation accuracy.

To clarify the assumptions underlying our formulation, we introduce a graphical model that represents all relevant random variables and their dependencies. This visualization provides a concise overview of the probabilistic relationships of our framework. Before presenting the graphical model, we define the symbols used in the diagram for clarity and consistency.

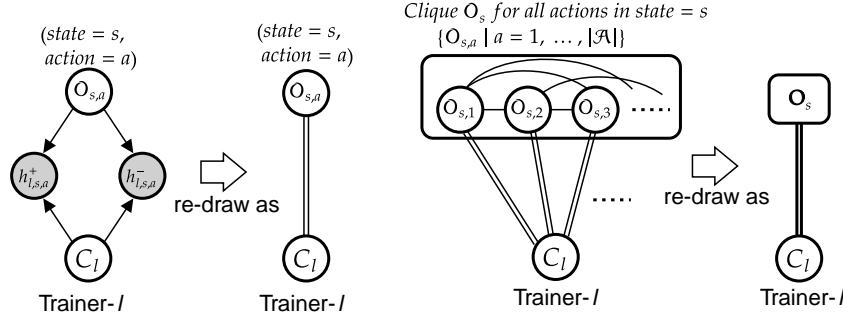


Figure 5: Symbols for graphical model for consistency level estimation with VI algorithm.

We assume that the number of positive and negative feedback for a given trainer and state-action pair ($h_{l,s,a}^+$ and $h_{l,s,a}^-$) are depending upon the optimality flag for the state-action pair ($O_{s,a}$) and the trainer's consistency level (C_l), and their relationship can be drawn as the left-most figure in Fig. 5. We redraw it with the double line between $O_{s,a}$ and C_l as the second left figure in Fig. 5. Next, we introduce a clique that has all optimality flags for a given state. The optimality flags in the clique are fully connected. Then, we rewrite the clique with the rounded box as shown in the rightmost figure in Fig. 5. Also, the connections between C_l and the optimality flags are replaced with a bold double line, as shown in the figure. Note that the bold double line has all related observed variables – the number of positive and negative feedback for all actions from a given trainer l in a state s .

With the above symbols, we can draw our graphical model for the random variables as Fig. 6.

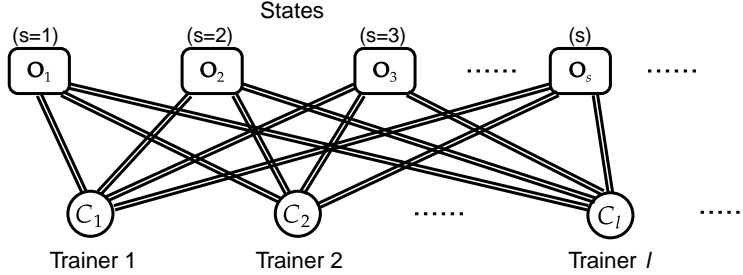


Figure 6: Graphical model for consistency level estimation with VI algorithm.

Based on the above graphical model, the complete likelihood function can be factorised as,

$$P(\mathbf{h}^+, \mathbf{h}^-, \mathbf{O}, \mathbf{C}) = \left(\prod_{s,a} \prod_l P(h_{l,s,a}^+, h_{l,s,a}^- | O_{s,a}, C_l) \right) \left(\prod_s P(\mathbf{O}_s) \right) \left(\prod_l P(C_l) \right), \quad (13)$$

The hidden variable \mathbf{O} is a matrix, and each element is an optimality flag $O_{s,a}$, which is a boolean that is 1 when a is the optimal action at state s , and 0 otherwise. The conditional probability of $h_{l,s,a}^+$ and $h_{l,s,a}^-$ is the binomial probability as defined

as follows:

$$P(h_{l,s,a}^+, h_{l,s,a}^- | \mathcal{O}_{s,a}; C_l) = \begin{cases} \begin{pmatrix} h_{l,s,a}^+ + h_{l,s,a}^- \\ h_{l,s,a}^+ \end{pmatrix} C_l^{h_{l,s,a}^+} (1 - C_l)^{h_{l,s,a}^-} & \text{if } \mathcal{O}_{s,a} = 1 \\ \begin{pmatrix} h_{l,s,a}^+ + h_{l,s,a}^- \\ h_{l,s,a}^- \end{pmatrix} (1 - C_l)^{h_{l,s,a}^+} C_l^{h_{l,s,a}^-} & \text{if } \mathcal{O}_{s,a} = 0, \end{cases} \quad (14)$$

and the prior of \mathbf{O}_s is set to the underlying RL algorithm policy ($P(\mathbf{O}_s = \mathbf{e}_a) = \pi_R(a|s)$). $P(C_l)$ is the prior for the l -th trainer's consistency level. We employ the conjugate prior (the beta distribution) for C_l , and it is defined as,

$$P(C_l) = \frac{\Gamma(\alpha_l + \beta_l)}{\Gamma(\alpha_l)\Gamma(\beta_l)} C_l^{\alpha_l - 1} (1 - C_l)^{\beta_l - 1}, \quad (15)$$

where Γ is the Gamma function. We approximate the posterior of \mathbf{O} and \mathbf{C} with two variational factors as follows, and VI performs approximate inference by updating each variational factor, $q(z)$, in turn, optimising the approximate posterior distribution until it converges.

$$P(\mathbf{O}, \mathbf{C} | \mathbf{h}^+, \mathbf{h}^-) = \left(\prod_s q(\mathbf{O}_s) \right) \left(\prod_l q(C_l) \right), \quad (16)$$

The variational factor $q(C_l)$ is updated by taking an expectation over the current estimate of the other variational factor $q(\mathbf{O}_s)$.

$$\log q(C_l) = (h_l^r + \alpha_l - 1) \log C_l + (h_l^w + \beta_l - 1) \log(1 - C_l) + \text{Constant}, \quad (17)$$

where the h_l^r and h_l^w are the expected number of right and wrong feedback, respectively. They are derived as,

$$\begin{aligned} h_l^r &= \sum_{s,a} q(\mathcal{O}_{s,a} = 1) h_{l,s,a}^+ + q(\mathcal{O}_{s,a} = 0) h_{l,s,a}^- \\ h_l^w &= \sum_{s,a} q(\mathcal{O}_{s,a} = 0) h_{l,s,a}^+ + q(\mathcal{O}_{s,a} = 1) h_{l,s,a}^- \end{aligned} \quad (18)$$

Constant in Eq. 17 is the partition function for $\log q(C_l)$. As it is a log of the beta distribution, The partition function becomes $\log \left(\frac{\Gamma(h_l^r + h_l^w + \alpha_l + \beta_l)}{\Gamma(h_l^r + \alpha_l)\Gamma(h_l^w + \beta_l)} \right)$. The other variational factor $q(\mathbf{O}_s)$ is updated as,

$$\log q(\mathbf{O}_s = \mathbf{e}_a) = \sum_l \left\{ \delta_{l,s,a} \mathbb{E} \log C_l - \delta_{l,s,a} \mathbb{E} \log(1 - C_l) \right\} + \log \pi_R(a|s) + \text{Constant}. \quad (19)$$

Where the $\mathbb{E} \log C_l$ and $\mathbb{E} \log(1 - C_l)$ are the expected log of the consistency level and one minus consistency level for the trainer l . They are derived as,

$$\begin{aligned} \mathbb{E} \log C_l &= \psi(h_l^r + \alpha_l) - \psi(h_l^r + h_l^w + \alpha_l + \beta_l) \\ \mathbb{E} \log(1 - C_l) &= \psi(h_l^w + \beta_l) - \psi(h_l^r + h_l^w + \alpha_l + \beta_l), \end{aligned} \quad (20)$$

where ψ is the digamma function. Constant in Eq. 19 is the partition function to make $\sum_a \log q(\mathbf{O}_s = \mathbf{e}_a) = 0$ for $\forall s$. Then, it repeats updating the variational factors until they converge. Once it is converged, we employ Eq. 19 to derive the posterior of \mathbf{O}_s and use it as the policy to decide the following action. The overall VI algorithm is summarised in Algorithm 1.

Appendix B: Derivation of Eq. 5

This appendix provides a detailed derivation of the posterior factorisation of the optimality variable, as expressed in Eq. 5. The derivation relies on the application of Bayes' rule in the first and third lines, as shown below. We also assume conditional independence between human feedback $\mathbf{h}^+, \mathbf{h}^-$ and the trajectories τ , given the optimality variable $\mathcal{O}_{s,a}$.

$$\begin{aligned} P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-, \tau) &= P(\mathbf{h}^+, \mathbf{h}^-, \tau | \mathcal{O}_{s,a}) \frac{P(\mathcal{O}_{s,a})}{P(\mathbf{h}^+, \mathbf{h}^-, \tau)} \\ &= P(\mathbf{h}^+, \mathbf{h}^- | \mathcal{O}_{s,a}) P(\tau | \mathcal{O}_{s,a}) \frac{P(\mathcal{O}_{s,a})}{P(\mathbf{h}^+, \mathbf{h}^-, \tau)} \quad (\because \text{the conditional independence}) \\ &= \frac{P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-) P(\mathbf{h}^+, \mathbf{h}^-)}{P(\mathcal{O}_{s,a})} \cdot \frac{P(\mathcal{O}_{s,a} | \tau) P(\tau)}{P(\mathcal{O}_{s,a})} \cdot \frac{P(\mathcal{O}_{s,a})}{P(\mathbf{h}^+, \mathbf{h}^-, \tau)} \\ &= P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-) P(\mathcal{O}_{s,a} | \tau) \frac{1}{P(\mathcal{O}_{s,a})} \frac{P(\mathbf{h}^+, \mathbf{h}^-) P(\tau)}{P(\mathbf{h}^+, \mathbf{h}^-, \tau)} \\ &\propto P(\mathcal{O}_{s,a} | \mathbf{h}^+, \mathbf{h}^-) P(\mathcal{O}_{s,a} | \tau) \frac{1}{P(\mathcal{O}_{s,a})}. \end{aligned} \quad (21)$$

Algorithm 1: Consistency Level Estimation with VI algorithm

Require: $i_{max}, \alpha, \beta, \pi_R(a|s), \mathbf{h}^+$ and \mathbf{h}^-

- 1: $\delta(l, s, a) \leftarrow h_{l,s,a}^+ - h_{l,s,a}^-$ for $\forall l, \forall s, \forall a$
- 2: $i \leftarrow 1$
- 3: Initialise $q(\mathbf{O}_s = \mathbf{e}_a) \leftarrow \pi_R(a|s)$ for $\forall s, \forall a$
- 4: **while** TRUE **do**
- 5: Update h_l^r and h_l^w for $\forall l$ with Eq. 18.
- 6: Update $\log q(C_l)$ for $\forall l$ with Eq. 17.
- 7: Update $\mathbb{E} \log C_l$ and $\mathbb{E} \log(1 - C_l)$ for $\forall l$ with Eq. 20.
- 8: Update $\log q(\mathbf{O}_s = \mathbf{e}_a)$ for $\forall s, \forall a$ with Eq. 19.
- 9: **if** converged($\mathbb{E} \log C_l, \mathbb{E} \log(1 - C_l), \forall l$) or $i == i_{max}$ **then**
- 10: break
- 11: **end if**
- 12: $i \leftarrow i + 1$
- 13: **end while**
- 14: **return** $\mathbb{E} \log C_l$ and $\mathbb{E} \log(1 - C_l)$ for $\forall l$

This final proportional expression highlights how the posterior probability of the optimality variable given both human feedback and trajectory observations can be decomposed into the product of two posterior terms, each conditioned on a different source of information, and adjusted by the prior probability of the optimality variable.

Appendix C: Evaluation Details

Here we describe our environmental setup for reproducibility of the experiments described in our paper. Across all environments, we compare our entropy-based active querying and crowd feedback aggregation approach against (1) the baseline approach (no feedback) and (2) random active learning (feedback randomly sampled from trajectories). We train an oracle policy as a proxy for an expert human trainer using standard reinforcement learning on the target environment. Full implementations are provided in our code appendix.

Gridworld Domains

We train both an (1) oracle as a synthetic human trainer and (2) an agent trained with the additional feedback generated from (1). Table 3 details our gridworld agent and oracle training parameters. Oracle training uses Greedy Tabular Q-Learning over a single trial.

PACMAN Our (5x5) PACMAN grid is illustrated in Figure 7a. Figure 7b shows feedback generation aggregated by location and compares these aggregations between randomly generated feedback and that generated by active learning. We see a strong concentration of feedback around the second piece of fruit and in the top left of the map.

Taxi Complementary evaluation featuring random starting positions and destinations and a larger map. Default parameters from the gymnasium library were used without modification Farama Gymnasium suite (Towers et al. 2024). Figure 8 shows the comparison of feedback gathered by random sampling and active learning. It appears that in this case (where random sampling and entropy-based sampling performed similarly) that entropy-based feedback seems to explore more. It could be that we are falsely attributing entropy earlier in learning when we have very little information. In the Taxi environment, individual states are less critical as the start and end position move between the red, blue, yellow and green tiles. Also there is no ghost chasing and forcing the agent into specific areas. This discrepancy in PACMAN and Taxi led us to investigate literature on critical states (Kumar et al. 2022; Goyal et al. 2023; Liu et al. 2023) and utilise Frozen Lake to create environments with varying numbers of bottlenecks. Further research is needed to align these theories and provide a comprehensive assessment.

Frozen Lake Customised versions of the Taxi environment from the Farama Gymnasium suite (Towers et al. 2024), designed to investigate the relationship between task difficulty and performance of our active learning method. The text version of our four maps are provided in Table 4.

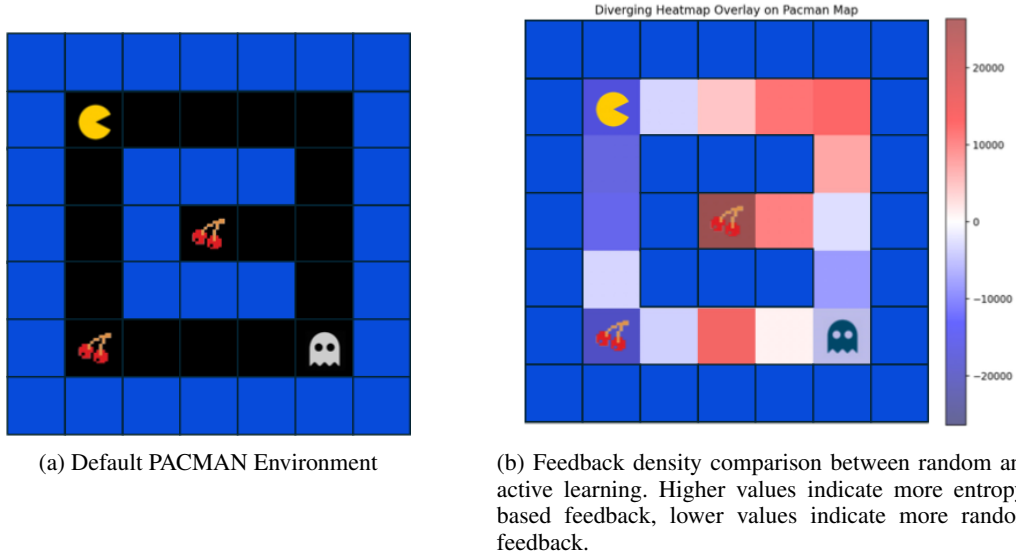


Figure 7: PACMAN Environment; (a) shows the starting positions of PACMAN and the ghost, as well as the two pellets which are to be collected to complete the game, (b) shows the different in feedback locations from random active learning and entropy-based active learning

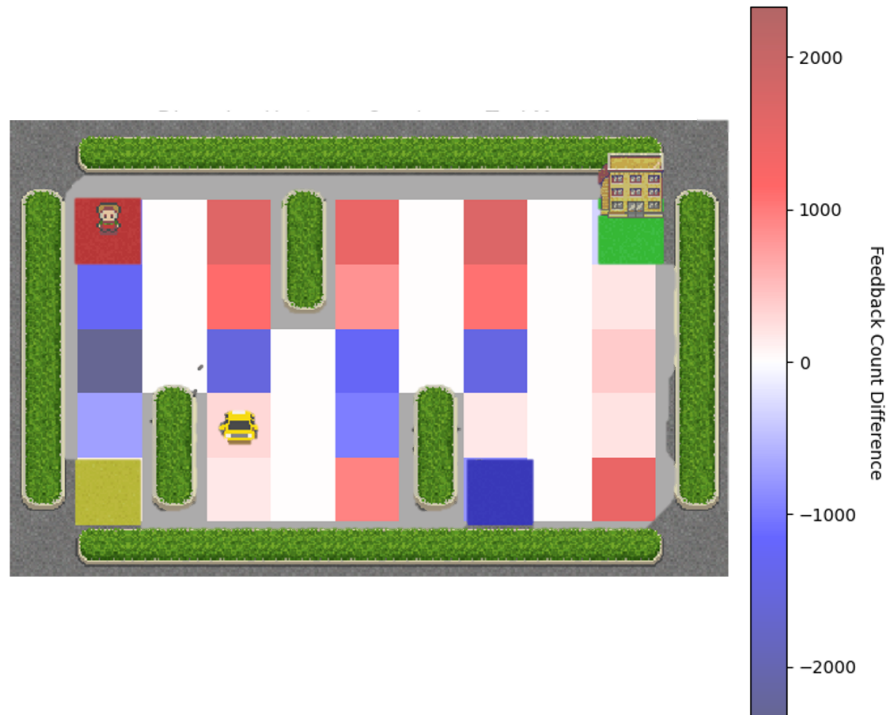


Figure 8: Taxi Feedback Comparison; comparing feedback sampled per location for random feedback selection and entropy-based feedback. The white ranks vertically adjacent to hedges should be ignored as they are not real states.

GridWorld Experimental Parameters	
Agent Training Parameters	
Learning Rate (α)	0.05
Discount Factor (γ)	0.9
Boltzmann Temperature (τ_b)	1.5
Prior of C_l	$\alpha_l = 90, \beta_l = 10$ for all trainers l
Number of Trainers	4
Trainer Consistency Levels (C_l)	[0.9, 0.8, 0.6, 0.3]
Episodes per Trial	1000
Number of Trials	50
Max Steps per Episode	500
Oracle Training Parameters	
Learning Rate (α)	0.05
Discount Factor (γ)	0.9
PACMAN	
Episodes	50000
Max Steps per Episode	500
Taxi	
Episodes	50000
Max Steps per Episode	500
Frozen Lake	
Episodes	5000
Max Steps per Episode	1000

Table 3: GridWorld experimental parameters for agent and oracle training.

(0) Random, Multiple Solutions	(1) Easy
SFFHFFFF	SFFHHHFF
FHFFFHFG	FFFFHHFG
FFHHFFFFH	FFFFFFFFF
FFFFFFHFF	FFFFHHFF
HFHFFHFF	FFFHHHFF
(2) Medium	(3) Hard (Maze)
SFHFFHHF	SFHFFFFH
FFHFFHFG	HFHHFHFG
FFHFFFFFF	HFHFFFFH
FFHFFHFF	HFHFHFHF
FFFFFFHHF	FFFFHFFF

Table 4: Frozen Lake Environment Variants

Type 1 Diabetes BGL Control

We use an UVA/Padova simulator (FDA-approved Type 1 Diabetes simulator) (Man et al. 2014), which takes meal and insulin injection information, then outputs a blood glucose level (BGL) as a CGM reading at each time step. The algorithms (agents) receive the meal, insulin and blood glucose level information and decide the amount of insulin taken in the next time step. We simulate the algorithms together with the Type 1 Diabetes simulator and evaluate how well the blood glucose levels are managed.

Environment Our simulator is based on an open-source implementation of UVA/Padova simulator (Xie. 2018), and we modified the state and action spaces as follows. For the action space, the insulin dose is discretised into six levels: 0, 20, 40, 80, 160, and 320 times the individual’s basal insulin dose.

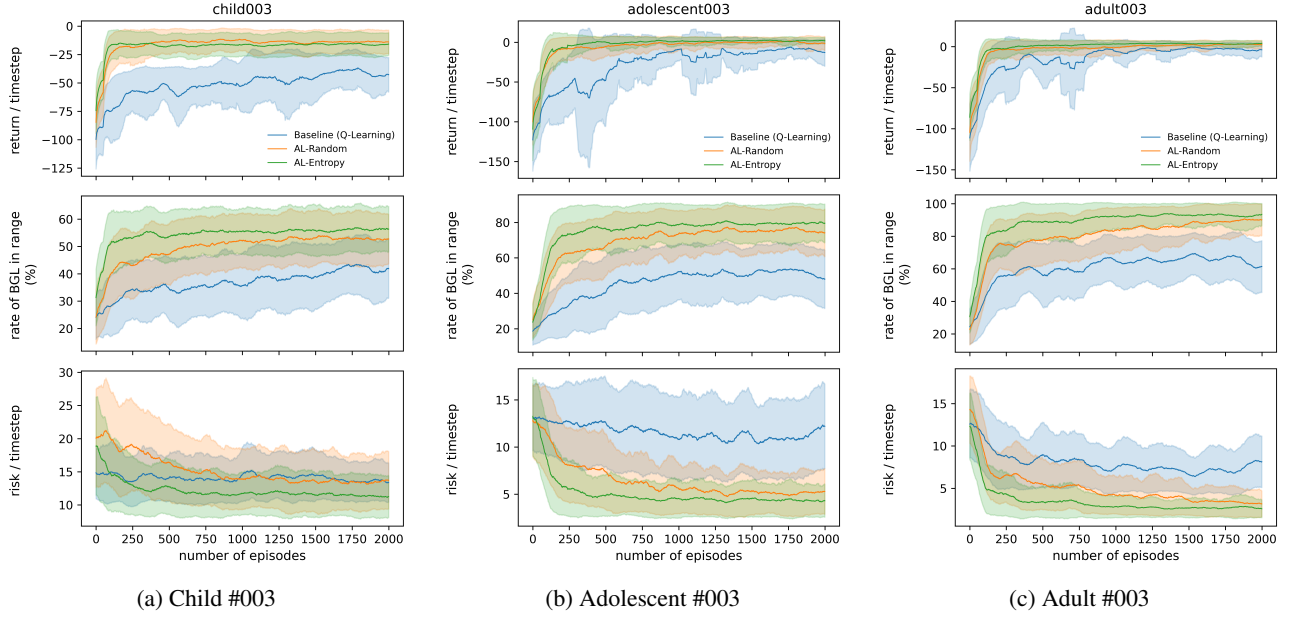


Figure 9: Simulation results for three person profiles with five learning trials for 2,000 episodes.

The state space is composed of the current BGL reading from CGM, the amount of carbohydrate intake, insulin on board (IOB) and time of day (ToD). Each of these elements is discretised as the following ranges:

- **BGL:** $[0, 70), [70, 90), [90, 110), [110, 180), [180, 300), [300, \infty]$ [mg/dL]
- **Carbs:** $[0, 2), [2, 10), [10, 20), [40, 60), [60, 80), [80, \infty]$ [g]
- **IOB:** $[0, 1), [1, 2), [2, 4), [4, \infty]$ [units]
- **ToD:** $[0, 6), [6, 12), [12, 18), [18, 24)$ [hour]

The IOB is the amount of active insulin remaining in the body, which is computed by 25% reduction of the amount of insulin every hour (assuming the insulin is active for four hours in the body). i.e. one insulin unit injected 1 hour ago and two units injected 3 hours ago:

$$\text{IOB} = 1 \times 0.75 + 2 \times 0.25 = 1.25$$

Experimental Setup We employ Q-Learning as the underlying RL algorithm combined with Boltzmann exploration policy. To simulate human feedback, we employ BBController as the oracle. The experiments are carried out 5 learning trials (each with different random seeds) for 2,000 episodes. We evaluate three virtual profiles: child#003, adolescent#003 and adult#003.

For this environment, we simulate 5 trainers with consistency levels of 0.9, 0.8, 0.7, 0.6, 0.3 who provide the feedback 10% of the time on average. Other hyper-parameters are set as follows:

- Learning Rate: $\alpha = 0.05$
- Discount Factor: $\gamma = 0.98$
- Boltzmann Exploration Temperature: $\tau_b = 10.0$
- Prior of C_l : $\alpha_l = 90, \beta_l = 10$ for all trainers l
- Base Variance of Q estimation: $\sigma_{base}^2 = 5000$
- Max Steps per Episode: 480
- CGM Hardware: Dexcom
- Insulin Pump: Insulet

Supplemental Results Figure 9 presents the simulation results for three virtual profiles: child#003, adolescent#003 and adult#003. Each simulation comprises five learning trials over 2,000 episodes. The figure is organised into three rows, presenting the following from top to bottom: the return (sum of reward) per time step, the percentage of time that BGL stays in the target range (70 - 180 mg/dL) and the averaged risk index per time step. The x-axis for all plots represents the number of episodes.

The figures show that our approach (AL-Entropy) outperforms other baselines (Baseline and AL-Random), which indicates the effectiveness of our entropy-based active feedback algorithm.